# ASCIWhiteGPFSstatus: August2001

*Richard Hedges, William Loewe,Tyce Mclarty, David Fox, Mark Grondona,Robin Goldstone, Terry Heidelberg*

**U.S. Department of Energy**

Lawrence
Livermore
National
Laboratory

**August 14, 2001**

## DISCLAIMER

This document was prepared as an account of work sponsored by an agency of the United States Government.  Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes.

This work was performed under the auspices of the U. S. Department of Energy by the University of California, Lawrence Livermore National Laboratory under Contract No. W-7405-Eng-48.

**ASCIWhiteGPFSstatus:August2001**

ThisreportbytheScalableI/OProject(SIOP)coverstheprogressofparallelI/O
performancefortheGPFSparallel            filesystemontheASCIWhitecomputer
system.

ThemajorhardwareandsoftwareconfigurationchangeswhichaffectGPFS
performanceare:(1)TheColonyswitchadaptershavebeenupgradedfrom
single/singletodouble/single;and(2)thesystemsoftwareha            sbeenupgraded
fromMohonk(PSSP3.2)withGPFS1.3usedduringFebruary2001testingtothe
presentMohonk2(PSSP3.3)withGPFS1.4.

Thetestsdiscussedherewereperformedonwhiteusingupto300computenodes
andtheGPFSfilesystemsusing16dedic            atedI/Onodes(servers).

ThetestsperformedutilizethePOSIXinterfacetoGPFS.Thenotedsystem
changestowhitehaveimprovedbothGPFSpeakreadandwritebyaboutafactor
oftwoforthelargescale(>=128nodes)tests.

**Machinecharacteristics**

AsofFebruary2001,whitenodeswereconnectedtotheColonyswitchby
single/singleadapters.By    lateJuly2001,thoseadaptershadbeenupgradedto
double/singleswitchadaptersandthesupportingsoftware(PSSP3.3andGPFS
1.4)wasinstalled.

White    has512totalnodesincluding16dedicatedGPFSI/Oservernodes.The
nodesare16  -waySMPnodescomprisedof375MhzPower3processors.Eachof
the16SSAservershas6RIOdrawerswith4SSAadaptersperdrawer.Eachof
these24adaptersisconnectedtoa            nSSAloop,whichconsistsof32disksspread
acrosstwoSSAenclosureswhereeachloopcontains64+Parraysand2hot
spares.[SeeFigure1.]Anindividualserverisresponsibleforserving3ofthesix
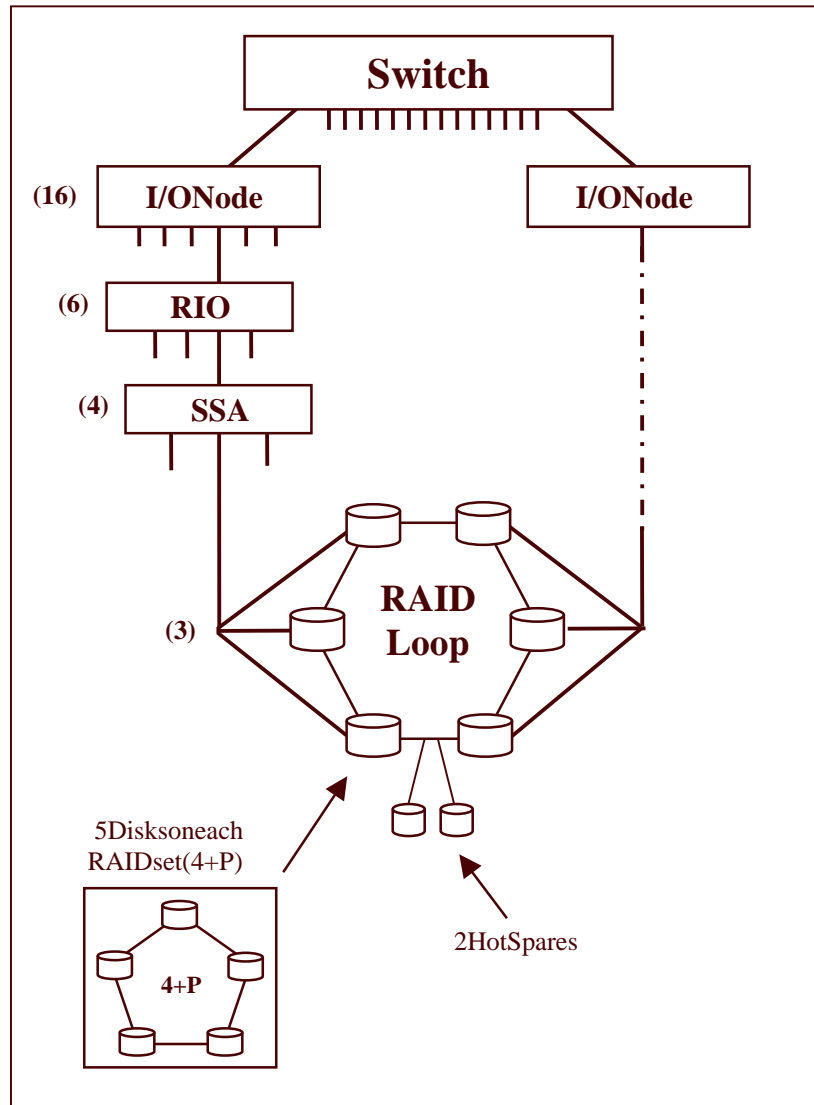disksinaloop,thusthereareatotalof1152dis            ksavailable:

> (3disks/loop*24loops/node*16nodes=1152disks)

Duetofilesystemlimitations,however,the/p/gw1filesystemactuallyconsisted
ofonly1024disksduringtheFebruarytesting.Toachievethisreductionindisks,
1diskperserver      wasremovedfromthe/p/gw1filesystemforeverythirdSSA
loop.Thisresultedintheremovalof8disksperserver,or128diskstotal,from
thefilesystemforafinalsizeof1024disks.

DuringthefirstsetoftestsonFebruary17,the/p/gw1files            ystembeingtested
consistedofonly1024disks.AsofJuly28,itconsistsofall1152disks.The

impactofthisconfigurationdifferenceisassumedtobeminimal:Thereismore
thanadequatediskbandwidthbeyondthethroughputcapacityoftheswitch.

Figure1diagramsthehardwarelayoutforthe/p/gw1filesystem.Thenumbersin
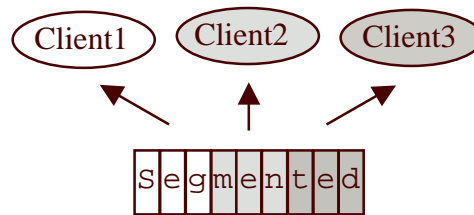parenthesestotheleftinthefigureshowthenumberofcomponentsateachlevel.



[Fig.1] *FileSystemHardware*

## TestDescription

ThesetestsusethePOSIXversionoftheIORtestcodewithsegmenteddata
patternandtransfersize=GPFSblocksize.Thisisthemostbasicandfamiliar
testusedbytheSIOPtoprobepeakI/OperformanceoftheGPF        Sfilesystemson

theASCImachines.ThepatternissegmentedwheneachprocessperformsI/O from/toalargecontiguousareaofthefiledistinctfromotherprocesses.[See Figure2.]



[Fig.2] *SegmentedPattern*

Moreexplicitlyw euse:segmentedpattern;1clientprocesspernode;varying nodes;create,writeandreadasinglecommonfile(size=512 $MB*n$).Each processaccesses512 $MB$,whicharecontiguousinthefilewhereforeachprocess $p$in0...n -1, $p$'sfirstbyteislocateda tlocation $p*512MB$.Thesizeofeach individualtransferistheGPFSblocksize=512 $KB$.


**Results**

February17,2001:single/single,IP

| Nodes | 30 | 60 | 128 | 256 |
|---|---|---|---|---|
| MB/swrite | 2473 | 3061 | 2170 | 2653 |
| MB/sread | 2589 | 3285 | 2306 | 3007 |


July28,2001: double/single,KLAPI

| Nodes | 32 | 64 | 128 | 300 |
|---|---|---|---|---|
| MB/swrite | 4262 | 5202 | 4793 | 5094 |
| MB/sread | 4862 | 4779 | 6921 | 6967 |


[Note:DifferentnodecountsforFebruaryandJulytests.]

University of California
Lawrence Livermore National Laboratory
Technical Information Department
Livermore, CA 94551